

# SciSpot: A Framework for Low-cost Scientific Computing On Transient Cloud Servers

Paper # 105

## ABSTRACT

In this paper, we...

## 1 INTRODUCTION

Scientific computing applications are a crucial component in the advancement of science and engineering, and play an important role in the analysis and processing of data, and understanding and modeling natural processes. These applications are typically implemented as large-scale parallel programs that use parallel-computing communication and coordination frameworks such as MPI. To take advantage of their parallel nature, conventionally, these applications have mostly been deployed on large, dedicated high performance computing infrastructure such as super computers.

(Vikram: Application of computer simulation plays a critical role in understanding natural and synthetic phenomena associated with a wide range of material,

biological, and engineering systems. This scientific computing approach involves the analysis and processing of data generated by the mathematical model representations of these systems implemented on computers. Typically, these scientific computing applications (simulations) are designed as parallel programs that leverage the communication and coordination frameworks associated with parallel computing techniques such as MPI in order to yield useful information at a faster pace (shorter user time). To exploit the parallel processing capabilities, such applications are routinely deployed on large, dedicated high performance computing infrastructure such as supercomputers [? ? ? ]. )

Increasingly, cloud computing platforms have begun to supplement *cite* and complement *cite; how true is this?* conventional HPC infrastructure ~~in-order~~ to meet the large computing and storage requirements of scientific applications. Public cloud platforms such as Amazon's EC2, Google Cloud Platform, and Microsoft Azure, offer multiple benefits such as *on-demand* resource allocation, convenient pay-as-you-go pricing models, ease of provisioning and deployment, and near-instantaneous elastic scaling. Most cloud platforms offer *Infrastructure as a Service*, and provide computing resources in the form of *Virtual Machines (VMs)*, on which a wide range of applications such as web-services, distributed data processing, distributed machine learning, etc., are deployed.

(Vikram: The extensive use of cloud platforms to host and run a wide range of applications such as web-services and distributed data processing have inspired early investigations in the direction of using these resources for scientific computing applications to meet the large computing and storage requirements of the latter. Early work in this area has shown the potential of using these cloud platforms to both supplement and complement conventional HPC infrastructure. Public cloud platforms such as Amazon's EC2, Google Cloud Platform, and Microsoft Azure, offer multiple benefits: *on-demand* resource allocation, convenient pay-as-you-go pricing models, ease of provisioning and deployment, and near-instantaneous elastic scaling, to name a few. Most cloud platforms offer *Infrastructure as a Service*, and provide computing resources in the form of Virtual Machines (VMs) that are application-agnostic and can serve as deployment sites for a wide range of applications

such as web-services, distributed machine learning, and scientific simulations.)

~~In order to~~ To meet the diverse resource demands of different applications, public clouds offer resources (i.e., VM's) with multiple different resource configurations (~~such as~~ e.g., number of CPU cores, ~~and~~ memory capacity), and pricing and availability contracts. Conventionally, cloud VMs have been offered with “on-demand” availability, such that the lifetime of the VM ~~wasis~~ solely determined by the owner of the VM (i.e., the cloud customer). Increasingly however, cloud providers have begun offering VMs with *transient*, rather than continuous on-demand availability. Transient VMs can be unilaterally revoked and preempted by the cloud provider, and applications running inside them face fail-stop failures. Due to their volatile nature, transient VMs are offered at steeply discounted rates. Amazon EC2 spot instances, Google Cloud Preemptible VMs, and Azure Batch VMs, are all examples of transient VMs, and are offered at discounts ranging from 50 to 90%.

However, deploying applications on cloud platforms presents multiple challenges due to the *fundamental* differences with conventional HPC clusters—which most applications still assume as their default execution environment. While the on-demand resource provisioning and pay-as-you-go pricing makes it easy to spin-up computing clusters in the cloud, ~~for effective resource utilization~~, the deployment of applications on cloud platforms must be cognizant of the heterogeneity in VM sizes, pricing, and availability for effective resource utilization. Crucially, optimizing for *cost* in addition to, ~~and not just~~ makespan, becomes an important objective in cloud deployments. Furthermore, although using transient resources can drastically reduce computing costs, their preemptible nature results in frequent job failures. Preemptions can be mitigated with additional fault-tolerance mechanisms and policies [17? ], although they impose additional performance and deployment overheads *such as? needs one more sentence to clearly explain why these existing approaches are not addressing the need that we address in this paper below....*

(Vikram: While the on-demand resource provisioning and pay-as-you-go pricing makes it easy to spin-up computing clusters in the cloud, the deployment of applications on cloud platforms must be cognizant of the heterogeneity in VM sizes, pricing, and availability for effective resource utilization. Crucially, optimizing for *cost* in addition to makespan, becomes an important objective in cloud deployments. Furthermore, although using transient resources can drastically reduce computing costs, their preemptible nature results in frequent job failures. Preemptions can be mitigated with additional

fault-tolerance mechanisms and policies [17? ], although they impose additional performance and deployment overheads. *add one more sentence.* These considerations of cost, server configuration heterogeneity, and frequent job failures intrinsic to the system present multiple challenges in deploying applications on cloud platforms which are fundamentally different from those that appear in using HPC clusters as the execution environment for the scientific computing applications.)

In this paper, we develop principled approaches for deploying ~~and orchestrating parallel~~ scientific computing applications on the cloud ~~at low cost~~, and present SciSpot, a ~~system framework~~ for low-cost scientific computing on ~~cloud~~ transient cloud servers. Our policies for tackling the resource heterogeneity and transient availability of cloud VMs build on a key insight: most scientific ~~com-~~puting applications are deployed as a collection or “bag” of jobs. These bags of jobs represent multiple instantiations of the same computation with different parameters. For instance, each job may be running a (parallel) simulation with a set of ~~simulation input~~ parameters, and different jobs in the collection run the (same) simulation ~~on employing~~ a different set of input parameters. Collectively, a bag of jobs can be used to “sweep” or search across a multi-dimensional parameter space to discover ~~or narrow down the set of feasible and viable~~ viable and/or interesting parameters associated with the modeled natural or synthetic processes. A similar approach is adopted in the use of machine learning (ML) to enhance scientific computational methods, a rapidly emerging area of research, when a collection of jobs with independent parameter sets are launched to train ML models to predict simulation results and/or accelerate the simulation technique.

Prior approaches and systems for mitigating transiency and cloud heterogeneity have largely targeted individual instantiations of jobs [17, 19? , 20]. For a bag of jobs, it is not necessary, or sufficient, to execute an individual job in timely manner—instead, we could selectively restart failed jobs in order to complete the necessary, desired subset ~~a fraction~~ of jobs in a bag. Furthermore, treating the bag of jobs as a fundamental unit of computation allows us to select the “best” server configuration for a given application, by exploring different servers for initial jobs and running the remainder of the jobs on the optimal server configuration *is this optimal server found on the fly after the initial set of explorations, or this can be found separately?*

We show that optimizing across an entire bag of jobs and being cognizant of the relation between different jobs in a bag, can enable simple and powerful policies for optimizing cost, makespan, and ease of deployment. We

implement these policies as part of the SciSpot framework, and make the following contributions: *skipping the following contributions, will look at it after they are reordered and further refined*

- (1) In order to select the “right” VM from the plethora of choices offered by cloud providers, we develop a search-based server selection policy that minimizes the cost of running applications. Our search based policy selects a transient server type based on it’s cost, parallel speedup, and probability of preemption.
- (2) Since transient server preemptions can disrupt the execution of jobs, we present the *first* empirical model and analysis of transient server availability that is *not* rooted in classical and out-dated bidding models for EC2 spot instances that have been proposed thus far. Our empirical model allows us to predict expected running and costs of jobs of different types and durations.
- (3) We develop preemption-mitigation policies to minimize the overall makespan of bags of jobs, by taking into consideration the partial redundancy and relative “importance” of different jobs within a bag. Combined, our policies yield a cost saving of XXX% and a makespan reduction of XXX% compared to conventional cloud deployments, and a makespan reduction of XXX% compared to a conventional HPC supercomputer.
- (4) Finally, ease of use and extensibility are one of the “first principles” in the design of SciSpot, and we present the design and implementation of the system components and present case studies of how scientific applications such as molecular dynamics simulations can be easily deployed on transient cloud VMs.

## 2 BACKGROUND AND OVERVIEW

In this section, we give an overview of the characteristics and challenges of transient cloud computing; motivate the need for the bag of jobs abstraction in scientific computing workflows; and give an overview of our SciSpot system.

### 2.1 Transient Cloud Computing

Infrastructure as a service (IaaS) clouds such as Amazon EC2, Google Public Cloud, Microsoft Azure, etc., typically provide computational resources in the form of virtual machines (VMs), on which users can deploy their applications. Conventionally, these VMs are leased on an “on-demand” basis: cloud customers can start up a VM when needed, and the cloud platform provisions and runs these VMs until they are shut-down by the customer. Cloud workloads, and hence the utilization of

cloud platforms, shows large temporal variations. To satisfy user demand, cloud capacity is typically provisioned for the *peak* load, and thus the average utilization tends to be low, of the order of 25% [? ?].

To increase their overall utilization, large cloud operators have begun to offer their surplus resources as low-cost servers with *transient* availability, which can be preempted by the cloud operator at any time (after a small advance warning). These preemptible servers, such as Amazon Spot instances [?], Google Preemptible VMs [?], and Azure batch VMs [?], have become popular in recent years due to their discounted prices, which can be 7-10x lower than conventional non-preemptible servers.

However, effective use of transient servers is challenging for applications because of their uncertain availability [? ?]. Preemptions are akin to fail-stop failures, and result in loss of the application’s memory and disk state, leading to downtimes for interactive applications such as web services, and poor throughput for long-running batch-computing applications. Consequently, researchers have explored fault-tolerance techniques such as checkpointing [17, 20?] and resource management techniques [19] to ameliorate the effects of preemptions for a wide range of applications. However, the effect of preemptions is dependent on a combination of application resource and fault model, and mitigating preemptions for different applications remains an active research area [12].

### 2.2 Bag of Jobs in Scientific Computing

The typical workflow associated with most scientific computing applications, often involves evaluating a computational model across a wide range of physical and computational parameters. For instance, constructing and calibrating a molecular dynamics application (such as [13]), usually involves running a simulation with different physical parameters such as characteristic sizes and interaction potentials, as well as computational parameters such as simulation timesteps. Each of these parameters can take a wide range of values, resulting in a large number of combinations which must be evaluated by invoking the application multiple times (also known as a parameter sweep). Since each computational job explores a single combination of parameters, this results in executing a “bag of jobs”, with each job in the bag running the same application, but with possibly different parameters.

The bag of jobs execution model is pervasive in scientific computing and applicable in many contexts. In

addition to exploratory parameter sweeps, bags of jobs also result from running the application a large number of times to account for model or computational stochasticity, and can be used to obtain tighter confidence intervals. Increasingly, bags of jobs also arise in the emerging research that combines statistical machine learning (ML) techniques and scientific simulations [2, 4, 6, 7, 13–16, 18, 23]. For instance, large bags of jobs are run to provide the necessary training and testing data for learning statistical models such as neural networks that are then used to improve the efficacy of the simulations.

The bag of jobs execution model has multiple characteristics, that give rise to unique challenges and opportunities when deploying them on cloud transient servers. First, since bags of jobs require a large amount of computing resources, deploying them on the cloud can result in high overall costs, thus requiring policies for minimizing the cost and overall running time. Second, we observe that usually, there is no dependency between individual jobs in a bag, thus allowing increased flexibility in job scheduling. And last, treating entire bags of jobs as an execution unit, instead of individual jobs, can allow us to use partial redundancy between jobs and reduce the fault-tolerance overhead to mitigate transient server preemptions.

## 2.3 SciSpot Overview

Our system, SciSpot, is a general-purpose software framework for running scientific computing applications on low-cost cloud transient servers. It incorporates policies and mechanisms for generating, deploying, orchestrating, and monitoring bags of jobs on cloud servers. Specifically, it runs a bag of jobs defined by these parameters:

Bag of job = { $\mathcal{A}$ : Application to execute,  
 $N$ : Number of jobs,  
 $m$ : Minimum number of jobs to finish,  
 $\pi$ : Generator function for job parameters,  
 $\mathcal{R}$ : Computing resources per job}

SciSpot seeks to minimize the cost and running time of bags of jobs of scientific computing applications. SciSpot’s cost and time minimizing policies for running bags of jobs are based on empirical and analytical models of the cost and preemption dynamics of cloud transient servers, which we present in the next section.

## 3 PREEMPTION DYNAMICS OF TRANSIENT CLOUD SERVERS

In order to understand and improve the performance of applications running on transient cloud servers, we must understand the nature and dynamics of their preemptions. The preemption characteristics are governed by

the supply of surplus resources, the demand for cloud resources, and the resource allocation policies enforced by the cloud operator. Therefore, in this section, we present empirical and analytical models to help us understand the nature of preemptions.

(Vikram: To measure and improve the performance of applications running on transient cloud servers, it is critical to understand the nature and dynamics of their preemptions. The preemption characteristics are governed by the supply of surplus resources, the demand for cloud resources, and the resource allocation policies enforced by the cloud operator. In this section, we present empirical and analytical models that describe these characteristics and enable an intuitive understanding of the nature of preemptions. )

### 3.1 Price-based preemption models

Amazon’s EC2 spot instances were the original cloud transient servers. The preemptions of EC2 spot instances ~~isare based on~~ controlled by their *price*, which is dynamically adjusted based on the supply and demand of cloud resources. Spot prices are ~~determined~~ based on a continuous second-price auction, and if the spot price increases above a pre-specified maximum-price, then the server is preempted.

Thus, the time-series of these spot prices can be used for understanding preemption characteristics such as the frequency of preemptions and the “Mean Time To Failure” of the spot instances. Many research projects have used publicly available<sup>1</sup>historical spot prices to characterize and model spot instance preemptions [? ? ]. For example, past work has analyzed spot prices and shown that the MTTF’s of spot instances of different hardware configurations and geographical zones ranges from a few hours to a few days [? ? ].

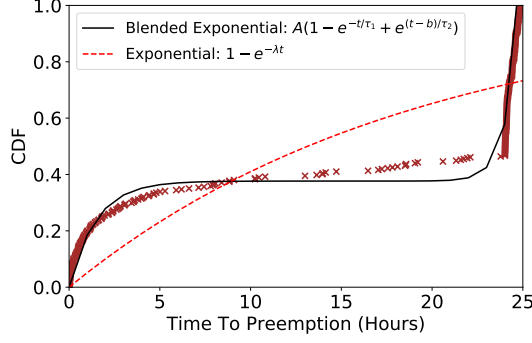
However, Amazon has recently changed the preemption characteristics of spot instances, and servers are now preempted even if the spot price is below the maximum price. Thus, spot prices are no longer a completely reliable indicator of preemptions, and preemptions can no longer be inferred from looking at prices alone. Therefore, new techniques are required to model preemption dynamics that can supplement the earlier price-based approaches, and we develop these techniques next.

### 3.2 Empirical models of preemptions

The preemptions of transient servers need not be related to their price. For example, Google’s Preemptible

<sup>1</sup>Amazon posts Spot prices of 3 months, and researchers have been collecting these prices since 2010 [? ? ].





**Figure 1: CDF of lifetimes of Google Preemptible Instances. Our blended exponential distribution fits much better than the conventional exponential failure distributions.**

VMs and Azure Batch VMs have a *fixed* price relative to their non-preemptible counterparts. In such cases, price based models are inadequate, and other approaches to understand preemptions *isare* required.

This task is further complicated by the fact that these cloud operators (Google and Microsoft) do not currently provide any information about preemption characteristics. Thus, relatively little is known about the preemptions (and hence the performance) of these transient VMs.

In order to understand preemption dynamics of transient servers, we conduct a large-scale empirical measurement study which is the first of its kind. We launched more than 1000 Google Preemptible VMs of different types over a two month period (Feb–April 2019), and measured their time to preemption (aka, their useful lifetime).<sup>2</sup>

A sample of 100 such preemption events are shown in Figure 1, which shows cumulative distribution of the VM lifetimes. Note that the cloud operator (Google) caps the *maximum* lifetime of the VM to 24 hours, and all the VMs are preempted before that hard limit. Furthermore, the lifetimes of VMs are *not* uniformly distributed, but have three distinct phases. *In the first phase, characterized by VM lifetime  $t \in (0, 3)$  hours, we observe that many VMs are quickly preempted after they are launched, and thus have a steep rate of failure initially; the rate of failure or preemptions can be obtained by taking the derivative of the CDF. The second phase characterizes the VMs that survive past 3 hours and enjoy a relatively low and uniform preemption rate over a relatively broad range of lifetime (characterized by the slowly rising CDF in*

Figure 1). The final phase exhibits a steep increase in the number of preemptions as the preemption deadline of 24 hours approaches. The overall rate of preemptions is “bath tub” shaped. *I think we should show the probability plot to exhibit the bath tub shape*

We note that this preemption behavior, imposed by the *constraint of the* small, 24 hour lifetime, is *substantially* different from conventional failure characteristics of hardware components and even EC2 spot instances. In these “classical” setups, the *rate-of-failure probability distribution characterizing failures* usually follows an exponential distribution  $f(t) = \lambda e^{-\lambda t}$ , where  $\lambda = 1/\text{MTTF}$ . Figure 1 shows the CDF ( $= 1 - e^{-\lambda t}$ ) of the exponential distribution when fitted to the observed preemption data, by finding the distribution parameter  $\lambda$  that minimizes the least squares error. From Figure 1, we can see that the classic exponential distribution is unable to model the observed preemption characteristics. *We attribute this deficiency to the central assumption made in the underlying reliability theory principles that leads to the exponential distribution: the rate of preemptions is independent of the lifetime of the VMs, in other words, the preemptions are memoryless.* This assumption breaks down when there is a fixed upper bound on the lifetime, as is the case for Google Preemptible VMs, and the conventional approach becomes insufficient to model this constrained preemption dynamics.

### 3.3 Analytical model of preemption dynamics in Google cloud

We now develop a *minimal* analytical model for preemption dynamics that is faithful to the empirically observed data and provides a basis for developing running-time and cost-minimizing optimizations presented in Section 4. This new model is based on the earlier observation that the cumulative distribution of lifetimes has multiple distinct temporal phases. The key assumption underlying our minimal model is the presence of two distinct failure processes that give rise to a new probability distribution characterizing the preemptions and the observed CDF, and ensure the dependence of the rate of failure on the VM lifetime. The first process dominates over the initial temporal phase and yields the classic exponential distribution that captures the steep rate of early preemptions. The second process dominates over the final phase near the 24 hour maximum VM lifetime and is assumed to be characterized by an exponential term that captures the sharp rise in preemptions that results from the constraint of a fixed 24 hour lifetime. Generally, these two processes compete during the middle phase to yield a relatively constant and low number of preemptions; in practice, based on the fits to the empirical data, we observe the

<sup>2</sup>We will release the complete preemption dataset and hope that other researchers can benefit.

first process to dominate over the second during this phase as well.

We propose the following general form for the CDF based on this model:

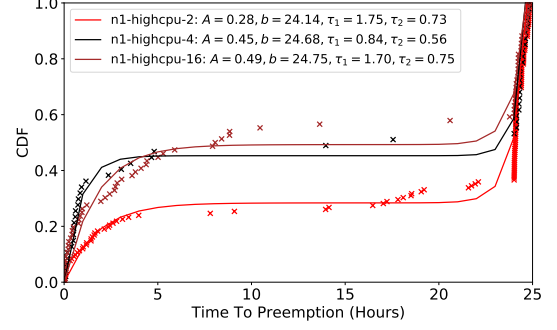
$$\mathcal{F}(t) = A \left( 1 - e^{-\frac{t}{\tau_1}} + e^{-\frac{t-b}{\tau_2}} \right), \quad (1)$$

where  $1/\tau_1$  is the rate of preemptions in the initial phase,  $1/\tau_2$  is the rate of preemptions in the final phase (generally,  $1/\tau_2 > 1/\tau_1$ ),  $b$  denotes the time when the preemptions occur at a high rate (generally, around 24 hours) which we term the activation time for the second process, and  $A$  is a constant used to scale the CDF to ensure that the initial conditions ( $F(0) = 0$ ) are met.

For most of its life, a VM sees failures according to the classic exponential distribution with a rate of failure equal to  $1/\tau_1$  – this behavior is captured by  $1 - e^{-t/\tau_1}$  term in Eq. 1. As VMs get closer to their maximum lifetime (24 hours) imposed by the cloud operator, they are reclaimed (i.e., preempted) at a high, exponential rate, which is captured by the second term introduced in the CDF ( $e^{-(t-b)/\tau_2}$ ). Shifting the argument ( $t$ ) of the exponential by  $b$  ensures that the exponential reclamation is only applicable towards the end of the VM’s maximum lifetime and does not dominate over the entire temporal range. As noted before,  $1/\tau_2$  is the rate of this reclamation.

The analytical model and the associated 4 parameter distribution function  $\mathcal{F}$  introduced above provides a much better fit to the empirical data and captures the different phases of the preemption dynamics through parameters  $\tau_1$ ,  $\tau_2$ ,  $b$ , and  $A$ . These parameters characterizing the preemption dynamics can be obtained for a given empirical CDF by minimizing least-squared function fitting methods.<sup>3</sup> In the next section, we use this analytical model for optimizing cloud resource selection such that we can run scientific applications at low cost and running times. We note that our motivation here is to provide a minimal model, i.e. a model based on data-driven observations and reasonable assumptions that provides a sufficiently accurate description of constrained preemption dynamics with the minimal number of necessary parameters. As is evident from Figure. 1, the analytical  $\mathcal{F}$  shows deviations from the data near the halfway point within the 24 hour lifetime. One can envision generalizing this model by including more failure processes characterized by failure rates and activation times (like  $b$ ) to capture the data with higher accuracy. Of course, this introduces a higher number of parameters

<sup>3</sup>More details about the distribution fitting are presented in the implementation section( ??



**Figure 2: The preemption characteristics of different VM types. Larger VMs are more likely to be preempted**

and reduces the predictive power and simplicity of the model.

### 3.4 Preemption dynamics of VMs of different types

Since cloud platforms support a wide range of applications, they also offer a large range of servers (VMs) with different resource configurations (such as the number of CPU cores, memory size, I/O bandwidths, etc.). For example, a cloud provider may offer VMs with (4 CPUs, 4 GB memory), (8 CPUs, 8 GB memory), etc. Most clouds offer a large number of different hardware configurations—Amazon EC2 offers more than 50 hardware configurations, for example [? ].

The expected lifetime of a preemptible VM is:

$$E[L] = \int_0^{24} t \mathcal{F}'(t) dt \quad (2)$$

Larger VM’s have a higher preemption probability and lower expected lifetimes. Because of [? ]

## 4 SCISPOT DESIGN

SciSpot handles all the cloud resource management and job scheduling associated with running a bag of jobs on transient cloud servers. In this section, we look at SciSpot’s policies for selecting the “right” cloud server for a given application, and policies for scheduling and running a bag of jobs on transient servers. Throughout, our aim is to minimize the overall cost and minimize the impact of preemptions.

SciSpot aims to provide a simple user interface to allow users to deploy their applications with the a minimum changes to their workflow. Most scientific applications are deployed on HPC clusters that have a cluster or a job manager such as Slurm [? ] or Torque [? ], and SciSpot integrates with the cluster manager (e.g., Slurm)

to provide the same interface to applications. As shown in Figure XXX TODO, SciSpot creates and manages clusters of transient cloud servers, and implements the various policies described in the rest of this section.

**High-level work flow:** When a user wishes to run a bag of jobs, SciSpot handles the provisioning of a cluster of transient cloud servers, and scheduling and monitoring of the bag of jobs, in addition to dealing with preemptions. Execution of a bag of jobs proceeds in two phases. In the first phase, SciSpot selects the “right” cluster configuration for a given application through a cost-minimizing exploration-based search policy, described next in Section 4.1. In the second phase, SciSpot proceeds to run the remaining jobs in the bag on the optimal cluster configuration.

## 4.1 Server Selection

**4.1.1 Why Server Selection is Necessary.** Before deploying any application on the transient cloud servers, we must first select the “right” cloud server for the application. Since cloud platforms support a wide range of applications, they also offer a large range of servers (VMs) with different resource configurations (such as the number of CPU cores, memory size, I/O bandwidths, etc.). For example, a cloud provider may offer VMs with (4 CPUs, 4 GB memory), (8 CPUs, 8 GB memory), etc. Most clouds offer a large number of different hardware configurations—Amazon EC2 offers more than 50 hardware configurations, for example [?].

*Importantly, different server configurations have different cost, performance, and preemption characteristics.*

Even if we assume that the total amount of resources to be allocated to a job is fixed, there are multiple *cluster configurations* to satisfy the allocation with the large number of available server types. For example, a job requiring a total of 128 CPUs can be run on a cluster of 2 servers with 64 CPUs each, or 4 servers with 32 CPUs each, etc. Server selection is especially important for parallel applications, because although the total amount of resources in each cluster configuration is constant, the resources are distributed differently. Since the performance of parallel applications is particularly sensitive to their communication overheads, different cluster configurations may yield different job running times. For instance, a smaller cluster with large servers will result in lower inter-server communication, and thus lower running times.

However, the performance of an application is also affected by preemptions of transient servers. Since preemptions are essentially fail-stop failures, synchronous parallel applications (such as those using MPI) are forced

to abort, and completing the job requires restarting it. Thus, preemptions can increase the overall running time of a job, with the increase in the job’s running time determined by the frequency of preemptions.

**4.1.2 Server Selection Policy.** Having provided the motivation and tradeoffs in server selection, we now describe how SciSpot’s server selection policy. Given an application and a bag of jobs, SciSpot “explores” and searches for the right server type by minimizing the expected cost of running the job.

We first determine the search space, which is the space of all cluster configurations  $\langle i, n_i \rangle$  such that  $r_i n_i = \mathcal{R}$ . With  $N$  server types, this results in at most  $N$  cluster configurations. Each configuration will yield different application performance, preemption overhead, and cost. Our server selection policy runs the application on on each configuration to determine its running time (in the absence of preemptions), which is denoted by  $T_{\langle i, n_i \rangle}$ . SciSpot thus does an exhaustive search over all valid cluster configurations to find the lowest-cost configuration  $\langle i, n_i \rangle$ .

We note that this search is different from conventional speedup plots in which the objective is to determine the how well an application scales with increasing amount of resources and parallelism. In contrast, we *fix* the total amount of resources allocated to the application’s job ( $= \mathcal{R}$ ), and only vary *how* these resources are distributed, which affects communication overhead and hence the performance. We assume that the total resource requirement for a job,  $\mathcal{R}$  can be easily provided by the user based on prior speedup data and the user’s cloud budget and the deadline for job completion.

To limit the search space, we observe that since most scientific applications are CPU bound, and we only need to consider VMs meant for CPU-bound workloads, such as **highcpu** VMs in Google Cloud and the **cc** family in Amazon EC2. For example, the Google cloud offers a total of 7 **highcpu** server types, yielding a small upper bound on the number of configurations to search.

**4.1.3 Server Cost Model.** Since server selection involves a tradeoff between cost, performance, and preemptions, we develop a model that allows us to optimize the resource allocation and pick the “best” server type that minimizes the expected cost of running an application on transient cloud servers.

Let us assume that the cloud provider offers  $N$  server types, with the price of a server type equal to  $c_i$ . Let  $\mathcal{R}$  denote the total amount of computing resources requested for the job. For ease of exposition, let us assume that  $\mathcal{R}$  is the total number of CPU cores. Furthermore,

let  $r_i$  denote the “size” of the server of type  $i$ . Then, the number of servers of type  $i$  required,  $n_i = \mathcal{R}/r_i$ .

The overall expected cost of running a job can then be expressed as follows:

$$E[C_{\langle i, n_i \rangle}] = n_i * c_i * E[T_{\langle i, n_i \rangle}] \quad (3)$$

Here,  $E[T_{\langle i, n_i \rangle}]$  denotes the expected running time of the job on  $n_i$  servers of type  $i$ . This running time, in turn depends on the preemption probability of the server type:

$$E[T_{\langle i, n_i \rangle}] = T_{\langle i, n_i \rangle} + P(\text{at least one failure}) * T_{\langle i, n_i \rangle} / 2 \quad (4)$$

Here,  $T_{\langle i, n_i \rangle}$  is the running time of the job without failures, which we obtain empirically as explained in the previous subsection.

The probability that at least one VM out of  $n_i$  will be preempted during the job execution can be expressed as:

$$P(\text{at least one failure}) = 1 - P(\text{no failure}) \quad (5)$$

$$= 1 - (1 - P(i, t))^{n_i} \quad (6)$$

$$= 1 - \left(1 - \frac{T_{\langle i, n_i \rangle}}{E[L_i]}\right)^{n_i} \quad (7)$$

Thus, we can see that if we select smaller VMs, we will require more of them (higher  $n_i$ ), and this cluster configuration will have a larger probability of failure and thus higher running times and costs.

Here,  $P(i, t)$  denotes the probability of a preemption of a VM of type  $i$  when a job of duration  $t$  runs on it.

It depends on the type of server, and we use historically determined failure distributions. In the expectation, it is given by:

$$P(i, t) = \frac{t}{E[L_i]} \quad (8)$$

Where  $E[L_i]$  is the expected lifetime of the VM of type  $i$  as computed with through the analytical model in earlier section. As a first order approximation, the running time of the job,  $t = T_{\langle i, n_i \rangle}$ , which is empirically obtained for a given application.

## 4.2 Scheduling a Bag of Jobs

In SciSpot’s

Once the right cluster configuration for a job has been determined, SciSpot then proceeds to run the remaining jobs in the bag. A bag of jobs is determined by the total number of jobs in the bag, associated parameters for each job, and the minimum number of jobs that must be successfully executed.

Given these parameters as input, SciSpot then

We note that the jobs in the bags are for the same application, so their computational and communication complexity can be assumed to the same across the bag.

## 5 SCISPOT INTERFACE AND IMPLEMENTATION

Central controller.

Cloud APIs for launching the jobs.

Slurm.

Job groups are specified via a JSON file that we then use to generate different parameter combinations.

Example of a JSON file here? What about the description? Maybe some details about ranges and fixed values?

## 6 EVALUATION

The contenders:

- (1) Run every job on un-tuned on-demand instance (cost and running time)
- (2) Run every job on nanohub/big-red-2 (running and waiting time)
- (3) Run on transient, restart every time (cost)
- (4) SciSpot with early stopping and job sacrificing

Performance of 3 benchmarks on different types of instance types and bigred2.

### 6.1 Preemption likelihood curves

### 6.2 Searching for the best cloud configuration

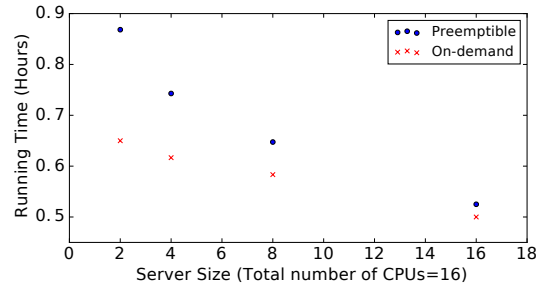


Figure 3: confinement running times

### 6.3 Total cost vs. running time graphs

### 6.4 HPC

## 7 RELATED WORK

### 7.1 Scientific applications on cloud

A classic survey is [11] [26]

Parameter sweep: [3]



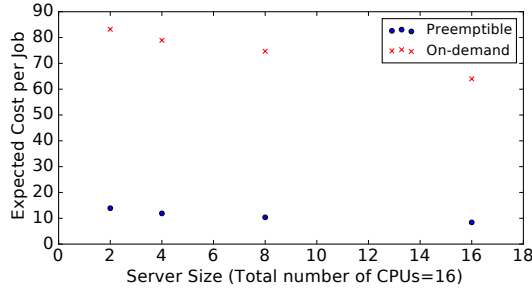


Figure 4: confinement cost

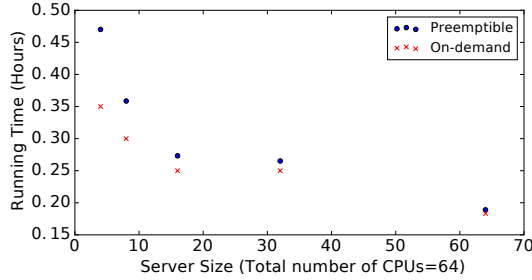


Figure 5: confinement running times

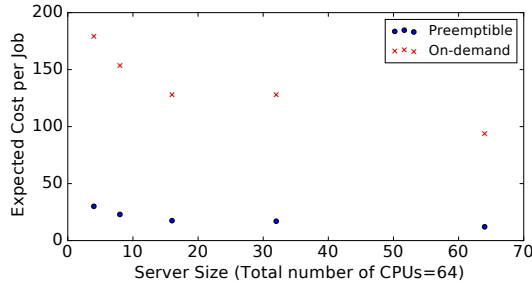


Figure 6: confinement cost

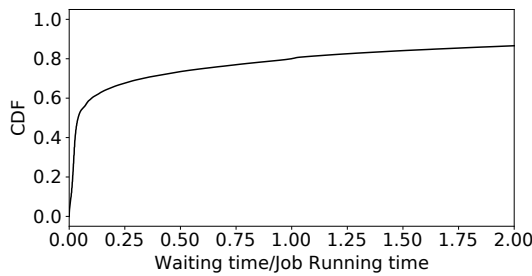


Figure 7: Ratio of waiting time to job running time on an HPC cluster. Average is 0.2

Price optimizations for Scientific workflows in the cloud [8]

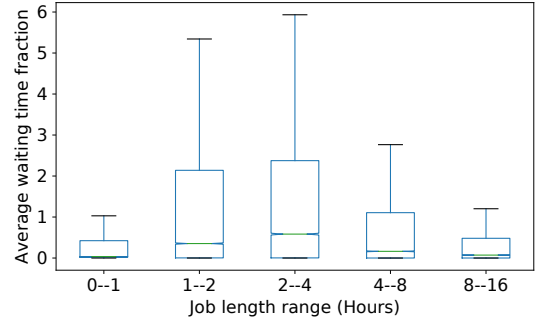


Figure 8: Waiting time fraction of jobs of different lengths varies.

## 7.2 Transiency mitigation

[17] classic work on MPI and Spot. Uses checkpointing. Redundancy, but for what? User specified number of VMs. Does not do instance selection. BCLR for checkpointing.

More spot and MPI: [9]. Focused on bidding and checkpoint interval. But bidding doesn't matter.

[21] is early work for spot and MPI and

[20] a batch computing service

Heterogeneity often used, but not useful in the context of MPI jobs [19]

Selecting the best instance type, often for data analysis computations [1], and [25], and others like Ernest and Hemingway.

All the past work was on EC2 spot market with gang failures and independent markets [9, 17]. However this assumption has now changed, and failures can happen anytime. Our failure model is more general, and applies to both cases.

**7.2.1 Fault-tolerance for MPI.** [5] has a discussion of checkpointing frequency which is comprehensive.

Replication is another way [22]

**7.2.2 Huge amount of work on bidding in HPC.** [24] [10]

## 7.3 Server Selection

Exploring a large configuration space using bayesian optimization methods in CherryPick [?] and Metis [?].

Can also use Latin Hypercube sampling for parameter exploration?

## REFERENCES

- [1] ALIPOURFARD, O., AND YU, M. CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics. 15.

- [2] BARTÓK, A. P., DE, S., POELKING, C., BERNSTEIN, N., KERMODE, J. R., CSÁNYI, G., AND CERIOTTI, M. Machine learning unifies the modeling of materials and molecules. *Science Advances* 3, 12 (2017).
- [3] CASANOVA, H., LEGRAND, A., ZAGORODNOV, D., AND BERMAN, F. Heuristics for scheduling parameter sweep applications in grid environments. In *Proceedings 9th Heterogeneous Computing Workshop (HCW 2000) (Cat. No.PR00556)* (May 2000), pp. 349–363.
- [4] CH'NG, K., CARRASQUILLA, J., MELKO, R. G., AND KHATAMI, E. Machine learning phases of strongly correlated fermions. *Phys. Rev. X* 7 (Aug 2017), 031038.
- [5] DONGARRA, J., HERAULT, T., AND ROBERT, Y. Fault tolerance techniques for high-performance computing. 66.
- [6] FERGUSON, A. L. Machine learning and data science in soft materials engineering. *Journal of Physics: Condensed Matter* 30, 4 (2017), 043002.
- [7] FOX, G., GLAZIER, J. A., KADUPITIYA, J., JADHAO, V., KIM, M., QIU, J., SLUKA, J. P., SOMOGYI, E., MARATHE, M., ADIGA, A., ET AL. Learning everywhere: Pervasive machine learning for effective high-performance computation. *arXiv preprint arXiv:1902.10810* (2019).
- [8] GAR, Y., MONGE, D. A., MATEOS, C., AND GARCA GARINO, C. Learning budget assignment policies for autoscaling scientific workflows in the cloud. *Cluster Computing* (Feb. 2019).
- [9] GONG, Y., HE, B., AND ZHOU, A. C. Monetary cost optimizations for MPI-based HPC applications on Amazon clouds: checkpoints and replicated execution. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis on - SC '15* (Austin, Texas, 2015), ACM Press, pp. 1–12.
- [10] GUO, W., CHEN, K., WU, Y., AND ZHENG, W. Bidding for Highly Available Services with Low Price in Spot Instance Market. In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing - HPDC '15* (Portland, Oregon, USA, 2015), ACM Press, pp. 191–202.
- [11] IOSUP, A., OSTERMANN, S., YIGITBASI, M. N., PRODAN, R., FAHRINGER, T., AND EPEMA, D. H. J. Performance Analysis of Cloud Computing Services for Many-Tasks Scientific Computing. *IEEE Transactions on Parallel and Distributed Systems* 22, 6 (June 2011), 931–945.
- [12] JOAQUIM, P., BRAVO, M., RODRIGUES, L., AND MATOS, M. Hourglass: Leveraging transient resources for time-constrained graph processing in the cloud. In *Proceedings of the Fourteenth EuroSys Conference 2019* (New York, NY, USA, 2019), EuroSys '19, ACM, pp. 35:1–35:16.
- [13] KADUPITIYA, J., FOX, G., AND JADHAO, V. *Submitted* (2018).
- [14] KADUPITIYA, J., FOX, G., AND JADHAO, V. Machine learning for performance enhancement of molecular dynamics simulations. Accepted.
- [15] LIU, J., QI, Y., MENG, Z. Y., AND FU, L. Self-learning monte carlo method. *Phys. Rev. B* 95 (Jan 2017), 041101.
- [16] LONG, A. W., ZHANG, J., GRANICK, S., AND FERGUSON, A. L. Machine learning assembly landscapes from particle tracking data. *Soft Matter* 11, 41 (2015), 8141–8153.
- [17] MARATHE, A., HARRIS, R., LOWENTHAL, D., DE SUPINSKI, B. R., ROUNTREE, B., AND SCHULZ, M. Exploiting redundancy for cost-effective, time-constrained execution of hpc applications on amazon ec2. In *HPDC* (2014), ACM.
- [18] SCHOENHOLZ, S. S. Combining machine learning and physics to understand glassy systems. *Journal of Physics: Conference Series* 1036, 1 (2018), 012021.
- [19] SHARMA, P., IRWIN, D., AND SHENOY, P. Portfolio-driven resource management for transient cloud servers. In *Proceedings of ACM Measurement and Analysis of Computer Systems* (June 2017), vol. 1, p. 23.
- [20] SUBRAMANYA, S., GUO, T., SHARMA, P., IRWIN, D., AND SHENOY, P. SpotOn: A Batch Computing Service for the Spot Market. In *SOCC* (August 2015).
- [21] TAIFI, M., SHI, J. Y., AND KHREISHAH, A. SpotMPI: A Framework for Auction-Based HPC Computing Using Amazon Spot Instances. In *Algorithms and Architectures for Parallel Processing*, Y. Xiang, A. Cuzzocrea, M. Hobbs, and W. Zhou, Eds., vol. 7017. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011, pp. 109–120.
- [22] WALTERS, J. P., AND CHAUDHARY, V. Replication-Based Fault Tolerance for MPI Applications. *IEEE Transactions on Parallel and Distributed Systems* 20, 7 (July 2009), 997–1010.
- [23] WARD, L., DUNN, A., FAGHANINIA, A., ZIMMERMANN, N. E., BAJAJ, S., WANG, Q., MONTROYA, J., CHEN, J., BYSTROM, K., DYLLA, M., ET AL. Matminer: An open source toolkit for materials data mining. *Computational Materials Science* 152 (2018), 60–69.
- [24] WOLSKI, R., BREVIK, J., CHARD, R., AND CHARD, K. Probabilistic guarantees of execution duration for Amazon spot instances. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis on - SC '17* (Denver, Colorado, 2017), ACM Press, pp. 1–11.
- [25] YADWADKAR, N. J., HARIHARAN, B., GONZALEZ, J. E., SMITH, B., AND KATZ, R. H. Selecting the best VM across multiple public clouds: a data-driven performance modeling approach. In *Proceedings of the 2017 Symposium on Cloud Computing - SoCC '17* (Santa Clara, California, 2017), ACM Press, pp. 452–465.
- [26] ZHAI, Y., LIU, M., ZHAI, J., MA, X., AND CHEN, W. Cloud versus in-house cluster: evaluating Amazon cluster compute instances for running MPI applications. In *State of the Practice Reports on - SC '11* (Seattle, Washington, 2011), ACM Press, p. 1.